

CUNY SPS DATA 621 - CTG5 - HW5

Betsy Rosalen, Gabrielle Bartomeo, Jeremy O'Brien, Lidiia Tronina, Rose Koh

May 15th, 2019

Contents

1	DATA EXPLORATION	2
1.1	Summary Statistics	2
1.2	Linearity	6
1.3	Missing Data	10
2	DATA PREPARATION	11
2.1	Missing Values	11
2.2	Transformation / Feature Engineering	11
3	BUILD MODELS	12
4	SELECT MODELS	13
4.1	Prediction	13
5	Appendix	14

Table 1: Data Dictionary

VARIABLE	DEFINITION	TYPE
TARGET	Number of Cases Purchased	count response
AcidIndex	Method of testing total acidity by using a weighted avg	continuous numerical predictor
Alcohol	Alcohol Content	continuous numerical predictor
Chlorides	Chloride content of wine	continuous numerical predictor
CitricAcid	Citric Acid Content	continuous numerical predictor
Density	Density of Wine	continuous numerical predictor
FixedAcidity	Fixed Acidity of Wine	continuous numerical predictor
FreeSulfurDioxide	Sulfur Dioxide content of wine	continuous numerical predictor
LabelAppeal	Marketing Score indicating the appeal of label design	categorical predictor
ResidualSugar	Residual Sugar of wine	continuous numerical predictor
STARS	Wine rating by a team of experts. 4 = Excellent, 1 = Poor	categorical predictor
Sulphates	Sulfate content of wine	continuous numerical predictor
TotalSulfurDioxide	Total Sulfur Dioxide of Wine	continuous numerical predictor
VolatileAcidity	Volatile Acid content of wine	continuous numerical predictor
pH	pH of wine	continuous numerical predictor

1 DATA EXPLORATION

When dining in a restaurant, a sommelier can assist you in selecting a perfect wine, even if you do not know much about wine yourself. By asking about your taste preferences, they can recommend a wine that pairs well with your meal, while complementing your likes and dislikes. But what happens when you're browsing the 12,000 commercially available wines as a large wine manufacturer, wondering how to select a good wine. For those who are not familiar with the industry, wine characteristics of the wine may only make selecting a product even more difficult.

The variables are mostly related to the chemical properties of the wine being sold, and for those who are not familiar with the industry, it may only make selecting a product even more difficult. The response variable is the number of sample cases of wine that were purchased by wine distribution companies after sampling a wine. These cases would be used to provide tasting samples to restaurants and wine stores around the United States.

1.1 Summary Statistics

Continuous and categorical variables were summarized separately for the sake of clarity.

STARS and LabelAppeal can be described as categorical variables. However, because of the numerical coding, these variables have been treated as if it were quantitative. It is best to designate such variables as factors so that they are treated appropriately.

Table 2: Summary statistics for numeric variables

	n	min	mean	median	max	sd
AcidIndex	12795	4.00	7.77	8.00	17.0	1.32
Alcohol	12142	-4.70	10.49	10.40	26.5	3.73
Density	12795	0.89	0.99	0.99	1.1	0.03
Sulphates	11585	-3.13	0.53	0.50	4.2	0.93
pH	12400	0.48	3.21	3.20	6.1	0.68
TotalSulfurDioxide	12113	-823.00	120.71	123.00	1057.0	231.91
FreeSulfurDioxide	12148	-555.00	30.85	30.00	623.0	148.71
Chlorides	12157	-1.17	0.05	0.05	1.4	0.32
ResidualSugar	12179	-127.80	5.42	3.90	141.2	33.75
CitricAcid	12795	-3.24	0.31	0.31	3.9	0.86
VolatileAcidity	12795	-2.79	0.32	0.28	3.7	0.78
FixedAcidity	12795	-18.10	7.08	6.90	34.4	6.32
TARGET	12795	0.00	3.03	3.00	8.0	1.93

Table 3: Summary statistics for categorical variables

	STARS	LabelAppeal
1 :3042	-2: 504	
2 :3570	-1:3136	
3 :2212	0 :5617	
4 : 612	1 :3048	
NA's:3359	2 : 490	

1.1.1 Summary Statistics Graphs

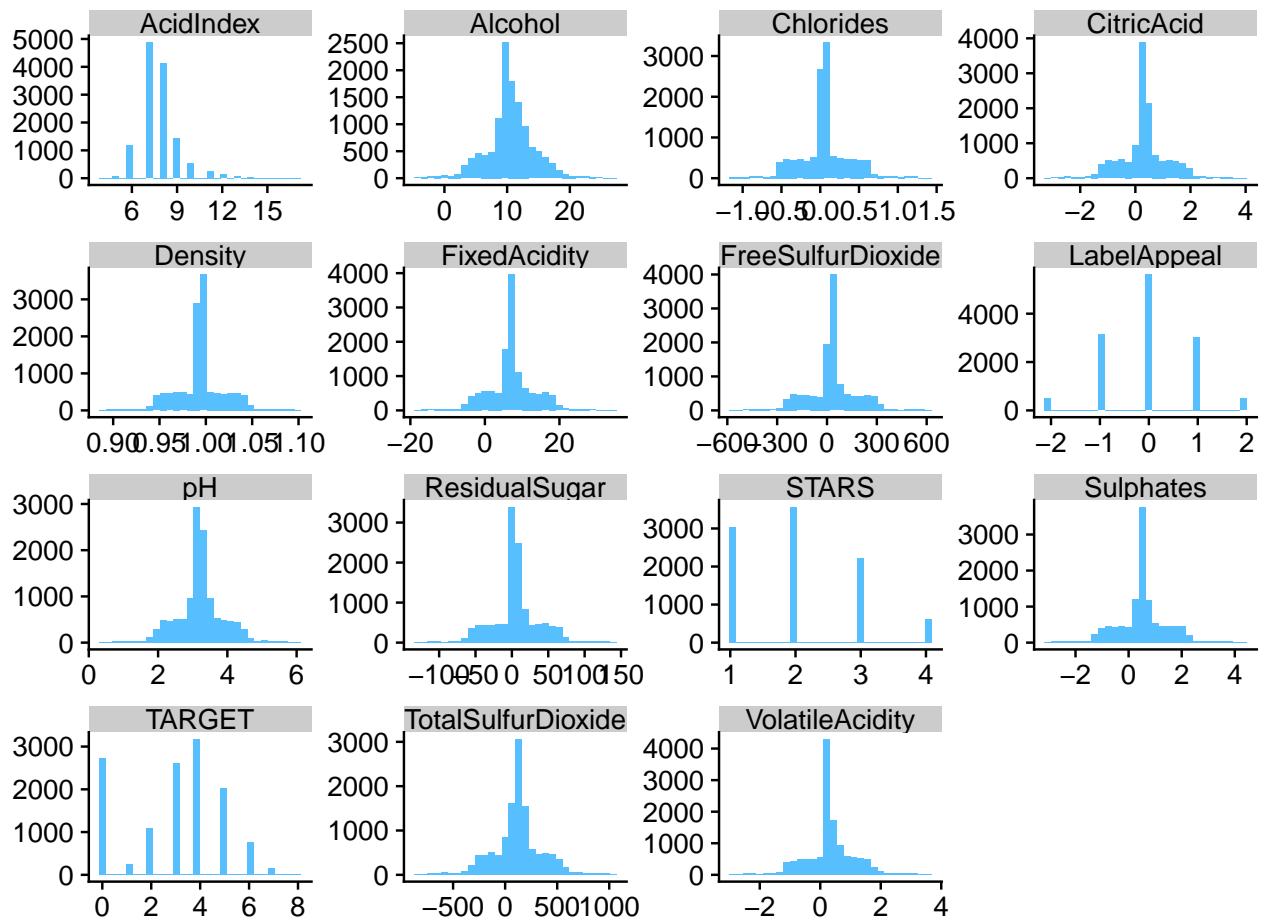


Figure 1: Data Distributions

Figure 1 shows the distribution of most of the variables seems normal although we see that the distribution is more peaked than a normal in Chlorides, CitricAcid, Density, FixedAcidity , FreeSulfurDioxide, ResidualSugar, Sulphates , ViolatileAcidity and AcidIndex have a slight right skew.

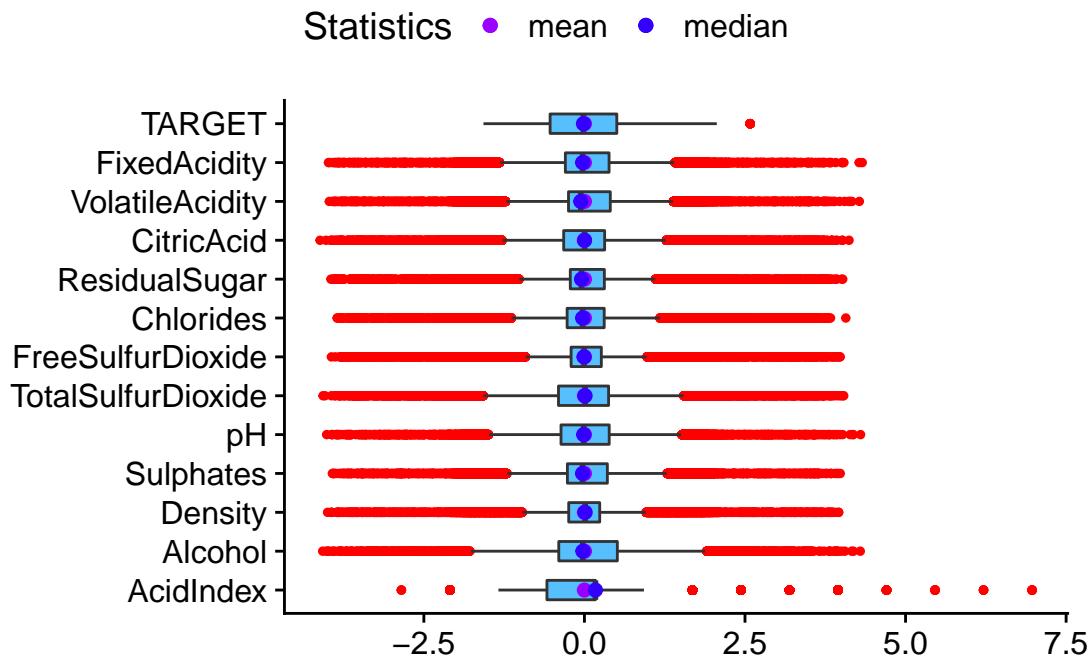


Figure 2: Scaled Boxplots

Figure 2 shows that there are a large number of outliers that need to be accounted for, except for `LabelAppeal`, `AcidIndex` and `STARS` which have limited number of variations.

1.2 Linearity

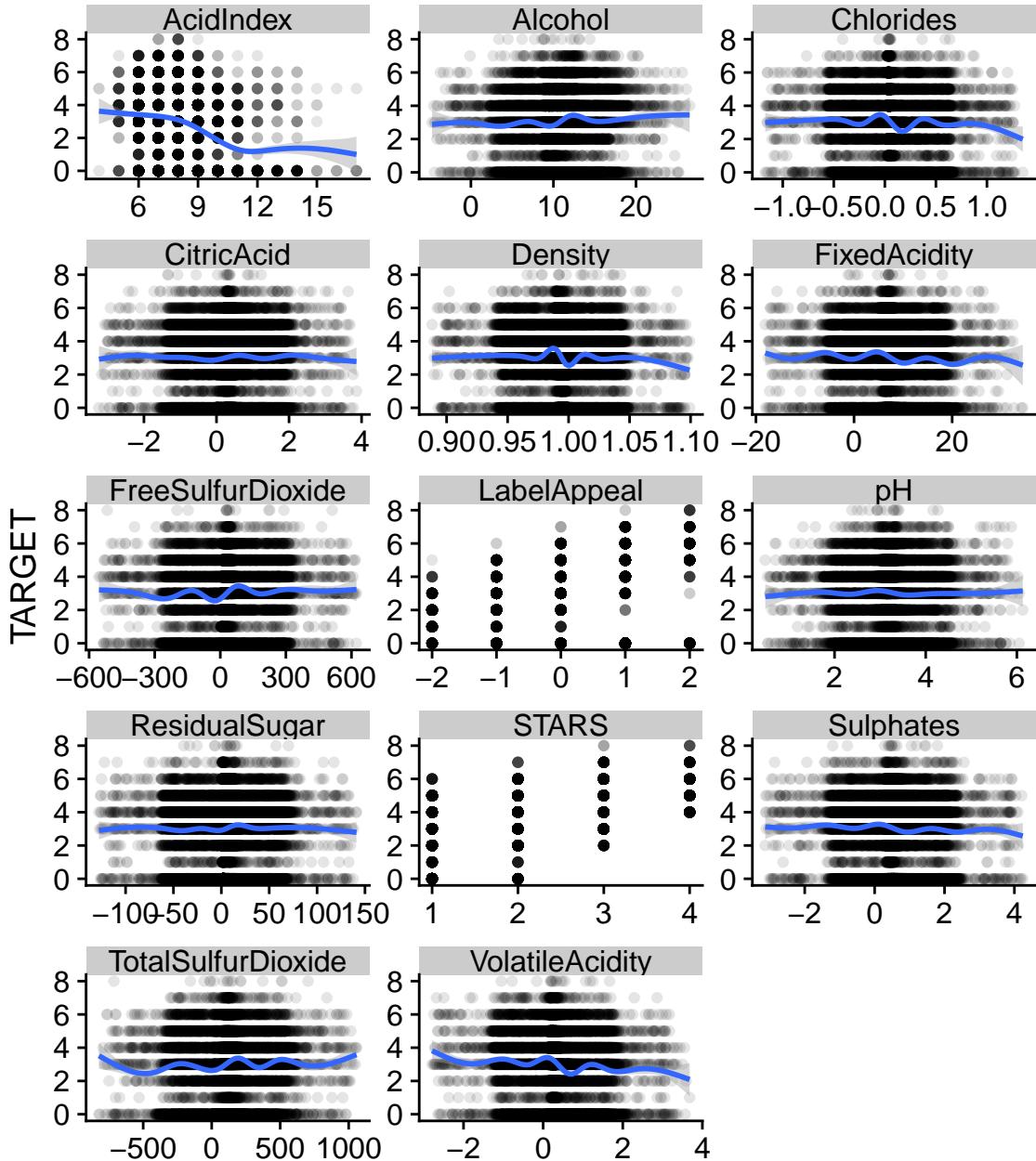


Figure 3: Scatter plot between numeric predictors and the TARGET

The raw predictors fail to show linear relationship with the TARGET except for the AcidIndex and VolatileAcidity. The Scatter Plots show a systematic, wave-like pattern for Density, FixedAcidity, FreeSulfurDioxide, TotalSulfurDioxide and VolatileAcidity.

1.2.1 Log Transformed Data

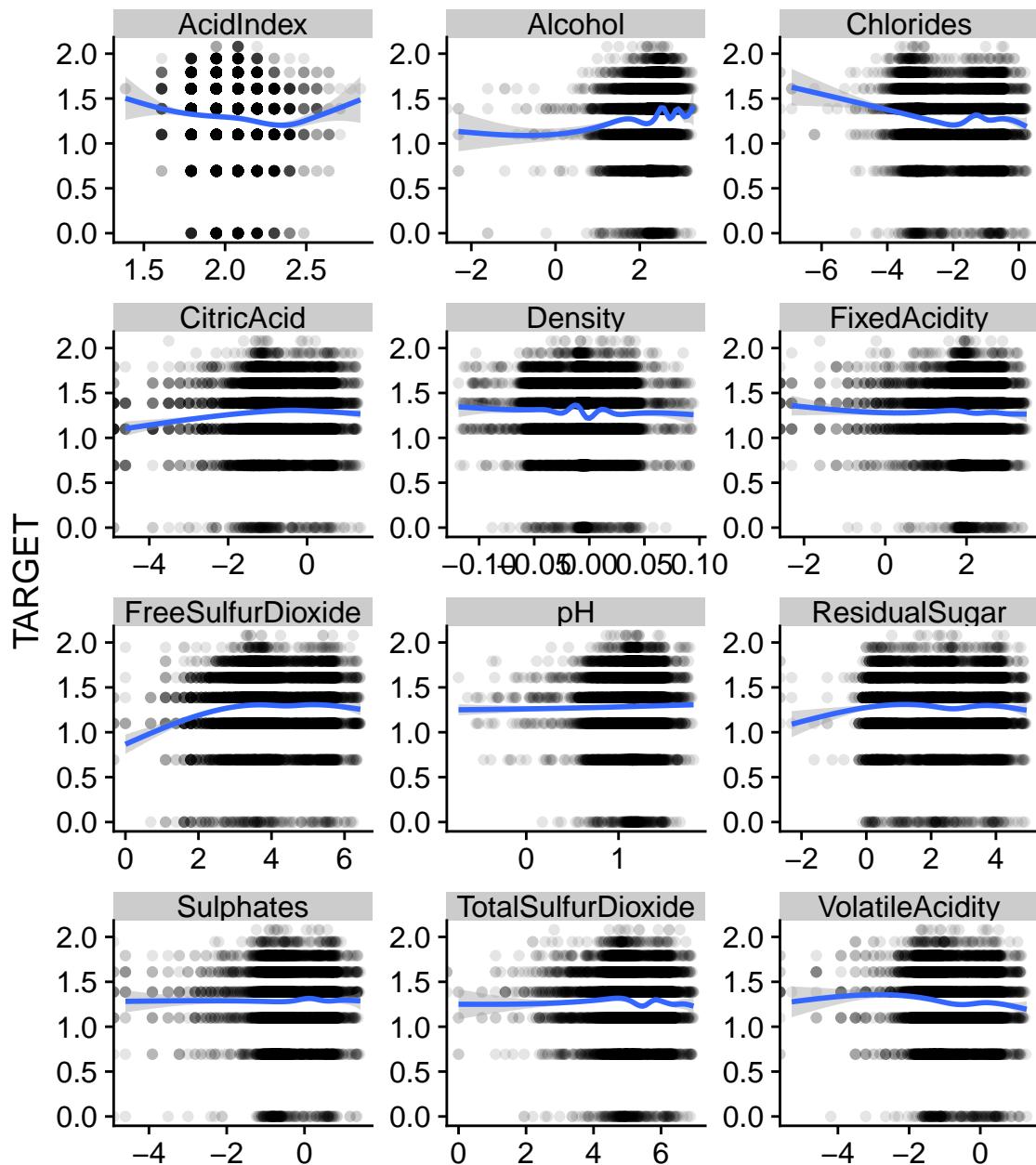


Figure 4: Scatter plot between log transformed predictors and the log transformed TARGET filtered for rows where TARGET is greater than 0

In attempt to improve the linearity of the variables against the TARGET variable, we start with a log transformation on all predictors and TARGET variable. As a result, the linearity of Chlorides and FreeSulfurDioxide become more apparent.

1.2.2 Box-Cox

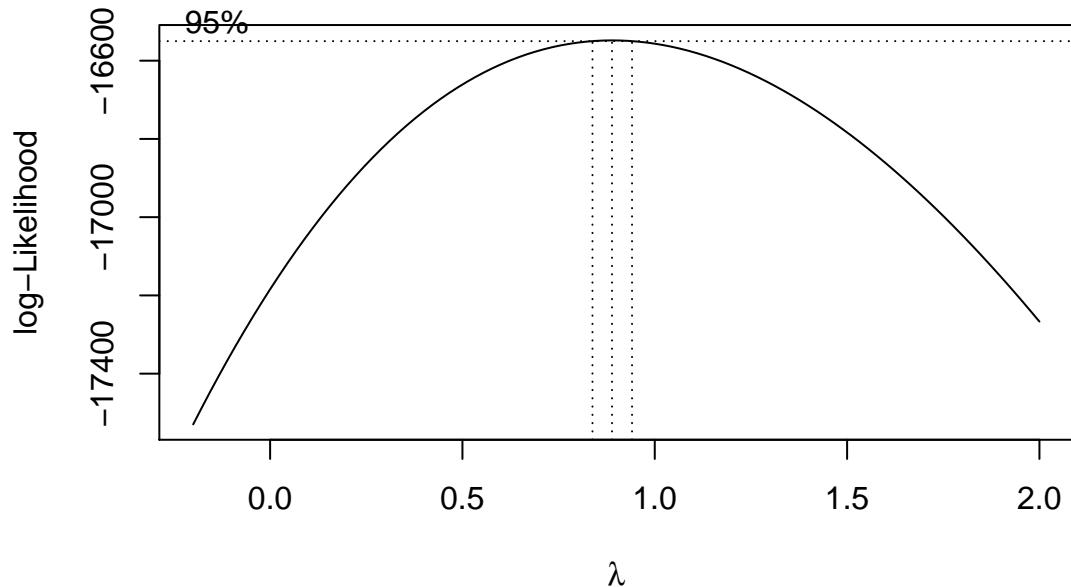


Figure 5: Box-Cox Plot

1.2.3 Square Root Transformed Predictors and Log Transformed Target

In ‘Linear Models with R’, Faraway suggested that the square root transformation is often appropriate for count response data. The Poisson distribution is a good model for counts, and that distribution has the property that the mean is equal to the variance thus suggesting the square root transformation. A plot of each predictor square root transformed plotted against the log transformed TARGET.

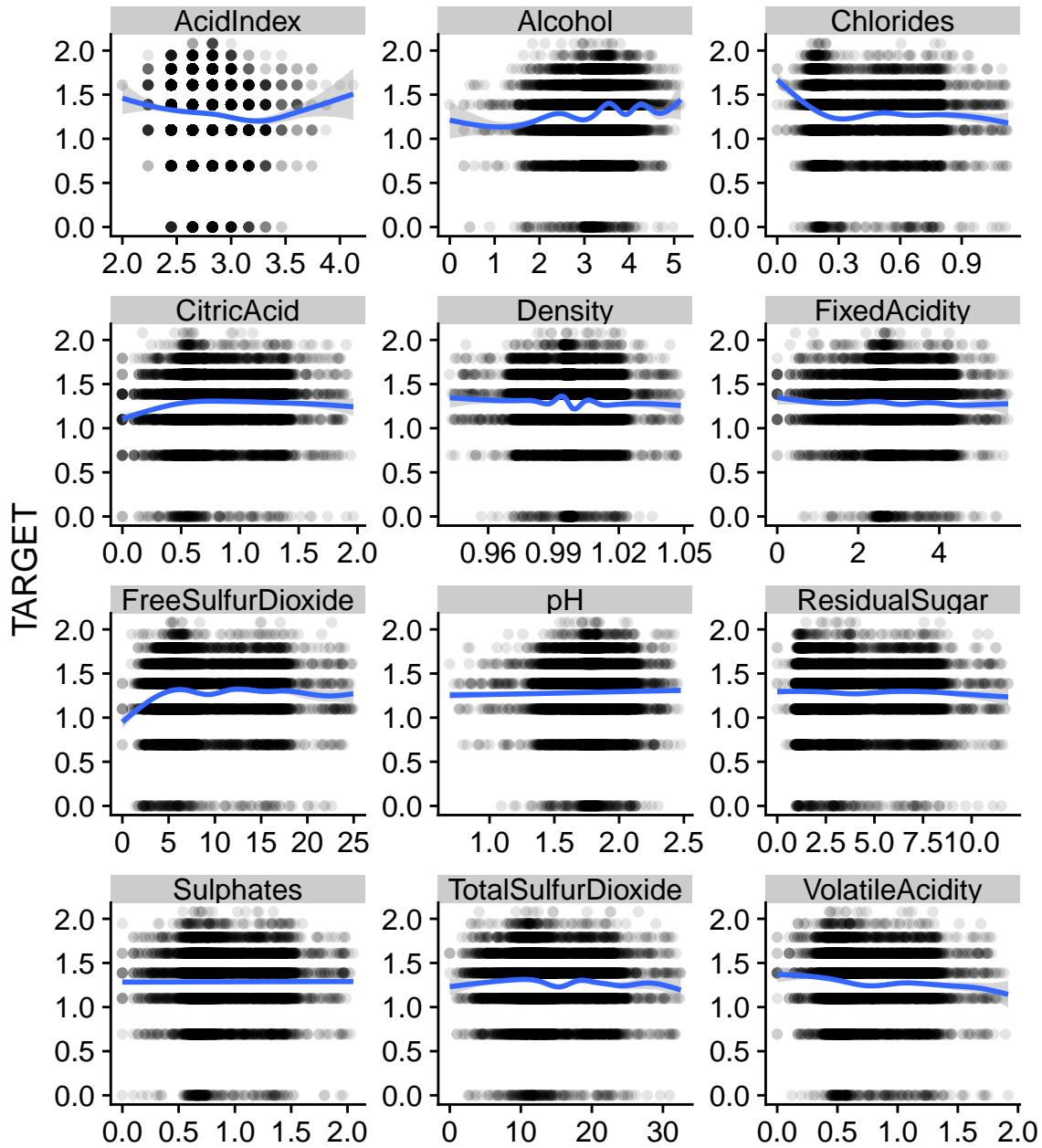


Figure 6: Scatter plot between square root transformed predictors and the square root transformed TARGET filtered for rows where TARGET is greater than 0

1.3 Missing Data

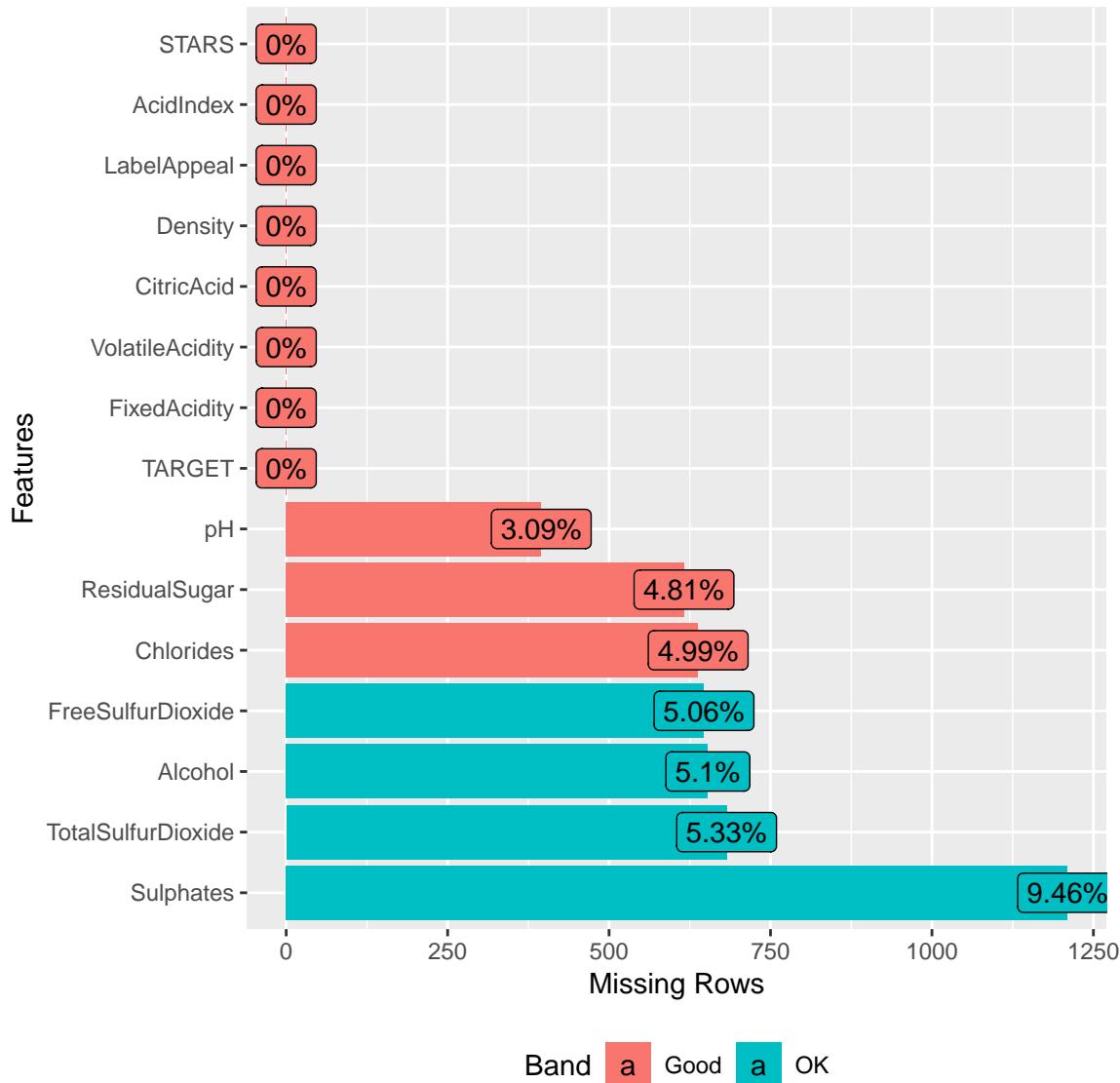


Figure 7: Missing data

A number of variables are missing observations: STARS, Sulphates, TotalSulfurDioxide, Alcohol, FreeSulfurDioxide, Chlorides, ResidualSugar, pH. For STARS, the number is 26.25%, but the others range between 3% and 9% of total. Approximately 50% of the cases are missing one of these variables.

2 DATA PREPARATION

Describe how you have transformed the data by changing the original variables or creating new variables

- a. Fix missing values (maybe with a Mean or Median value)
- b. Create flags to suggest if a variable was missing
- c. Transform data by putting it into buckets
- d. Mathematical transforms such as log or square root (or use Box-Cox)
- e. Combine variables (such as ratios or adding or multiplying) to create new variables

[JO: Consider strategies to transform negative variables - arithmetic?]

2.1 Missing Values

[Rough explanation] About half of the cases are missing at least one of the following variables: STARS, Sulphates, TotalSulfurDioxide, Alcohol, FreeSulfurDioxide, Chlorides, ResidualSugar, or pH. We use MICE (Multivariate Imputation by Chained Equations) to impute values these values based on ... [JO following up]

Pink and blue lines indicate close fit match in distribution of imputed values and recorded values, with the exception of STARS - the addition of STARS values has given the appearance of shifting the distribution from high to low lambda values.

2.2 Transformation / Feature Engineering

There are a large number of negative values for variables for which that is nonsensical, examples: Alcohol, CitricAcid, FixedAcidity, FreeSulfurDioxide, ResidualSugar, Sulphates, TotalSulfurDioxide. and VolatileAcidity. The range for the Poisson and negative binomial distribution has zero as a lower bound, so we can arithmetically transform the aforementioned variables to scale the lower IQR non-outlier values from zero up and drop the sub-IQR values (now the only negative values remaining). [JO currently function not working as intended - moving all points up by the mean - IQR * 1.5, and tossing anything less than that - so tuning calculations]

(<https://www.imachordata.com/do-not-log-transform-count-data-bitches/>)

Alternatively, we can explore whether more information on how measurements were made can be found to discern whether there might be some reason for negative values, and if there's a possibility of systematic data errors

We'll also test the data for overdispersion when setting up negative binomial models. [CODING UNDERWAY]

3 BUILD MODELS

Using the training data set, build at least two different poisson regression models, at least two differ-

Discuss the coefficients in the models, do they make sense? In this case, about the only thing you can c

4 SELECT MODELS

Decide on the criteria for selecting the best count regression model. Will you select models with slight

For the count regression model, will you use a metric such as AIC, average squared error, etc.? Be sure

4.1 Prediction

5 Appendix

The appendix is available as script.R file in `project5_wine` folder.

https://github.com/betsyrosalen/DATA_621_Business_Analyt_and_Data_Mining